

AD-A248 104



NEURAL COMPUTATION (in press)

①

Object Discrimination based on Depth-from-Occlusion

N00014-90-J-1864

Leif H. Finkel and Paul Sajda

1991

DTIC  
ELECTE  
APR 01 1992  
S D D

This document has been approved  
for public release and sale, its  
distribution is unlimited.

Department of Bioengineering and  
Institute of Neurological Sciences  
University of Pennsylvania  
Philadelphia, PA. 19104-6392

02 2 22 017

92-07824



We present a model of how objects can be visually discriminated based on the extraction of depth-from-occlusion. Object discrimination requires consideration of both the binding problem and the problem of segmentation. We propose that the visual system binds contours and surfaces by identifying "proto-objects"—compact regions bounded by closed contours. Proto-objects can then be linked into larger structures. The model is simulated by a system of interconnected neural networks. The networks have biologically-motivated architectures and utilize a distributed representation of depth. We present simulations that demonstrate three robust psychophysical properties of the system. The networks are able to stratify multiple occluding objects in a complex scene into separate depth planes. They bind the contours and surfaces of occluded objects (for example, if a tree branch partially occludes the moon, the two "half-moons" are bound into a single object). Finally, the model accounts for human perceptions of illusory contour stimuli.

## 1 Introduction

In order to discriminate objects in the visual world, the nervous system must solve two fundamental problems: binding and segmentation. The binding problem (Barlow, 1981) addresses how the attributes of an object—shape, color, motion, depth—are linked to create an individual object. Segmentation deals with the converse problem of how separate objects are distinguished. These two problems have been studied from the perspectives of both computational neuroscience (Marr, 1982; Grossberg and Mingolla, 1985; T. Poggio, et al., 1988; Finkel and Edelman, 1989) and machine vision (Guzman, 1968; Rosenfeld, 1988; Aloimonos and Shulman, 1989; Fisher, 1989). However, previous studies have not addressed what we consider to be the central issue: how does the visual system define an object—i.e., what constitutes a "thing".

Object discrimination occurs at an intermediate stage of the transformation between 2D image intensity values and visual recognition, and in general, depends upon cues from multiple visual modalities. In order to simplify the problem, we restrict ourselves to discrimination based solely on occlusion relationships. In a typical visual scene, multiple objects may occlude one another. When this occurs, it creates a perceptual dilemma—to which of the two overlapping surfaces does the common border belong? If the border is, in fact, an occlusion border, then it belongs to the occluding object. This identification results in a stratification of the two

1

Statement A per telecon  
Dr. Harold Hawkins ONR/Code 1142  
Arlington, VA 22217-5000

NWW 3/31/92



Dist		Avail and/or Special	
A-1			

30	
21	
1	

Identity Group

objects in depth and a de facto discrimination of the objects. Consider the case of a tree branch crossing the face of the moon. We perceive the branch as closer and the moon more distant, but in addition, the two "half-moons" are perceptually linked into one object. The visual system supplies a virtual representation of the occluded contours and surfaces in a process Kanizsa (1979) has called "amodal completion". With this example in mind, we propose that the visual system identifies "proto-objects" and determines which proto-objects, if any, should be linked into objects. For present purposes, a proto-object is defined as a compact region surrounded by a closed, piecewise continuous contour and located at a certain distance from the viewer. The contour can be closed upon itself, or more commonly, it can be closed by termination upon other contours.

We will demonstrate how a system of interconnected, physiologically-based neural networks can identify, link, and stratify proto-objects into objects. The networks operate, largely in parallel, to carry out the following interdependent processes:

- discriminate edges
- segment and bind contours
- identify proto-objects (i.e., bind contours and surfaces)
- identify possible occlusion boundaries
- stratify occluding objects into different depth planes
- attempt to link proto-objects into objects
- influence earlier steps (e.g. contour binding) by results of later steps (e.g. object linkage)

The constructed networks implement these processes using a relatively small number of neural mechanisms (such as detecting junctions of contours, and determining which surface is inside a closed contour). A few of the mechanisms used are similar to those of previous proposals (Grossberg and Mingolla, 1985; Finkel and Edelman, 1989; Fisher, 1989). But our particular choice of mechanisms is constrained by two considerations. First, we utilize a distributed representation of depth—this is based on the example of how disparity is represented in the visual

cortex (G. Poggio, et al, 1988; Lehky and Sejnowski, 1990). The relative depth of a particular object is represented by the relative activation of corresponding units in a *foreground* and *background* map. Second, as indicated above, we make extensive use of feedback (reentrant) connections from higher level networks to those at lower levels—this is particularly important in linking proto-objects. For example, once a higher level network has determined an occlusion relationship it can modify the way in which an earlier network binds contours to surfaces.

Any model of visual occlusion must be able to explain the perception of illusory (subjective) contours, since these illusions arise from artificially arranged cues to occlusion (Gregory, 1972). The proposed model can account for the majority of such illusions. In fact, the ability to link contours in the foreground and background corresponds, respectively, to the processes of modal and amodal completion hypothesized by Kanizsa (1979). The present proposal differs from previous neural models of illusory contour generation (Ullman, 1976; Grossberg and Mingolla, 1985; von der Heydt, et al., 1989; Finkel and Edelman, 1989) in that it generates illusory objects—not just the contours. The difference is critical: a network which generates responses to the three sides of the Kanizsa triangle, for example, is not representing a triangle (the object) *per se*. To represent the triangle it is necessary to link these three contours into a single entity, to know which side of the contour is the inside, to represent the surface of the triangle, to know something about the properties of the surface (its depth, color, texture, etc.), and finally to bind all these attributes into a whole. This is clearly a much more difficult problem. We will describe, however, a simple model for how such a process might be carried out by a set of interconnected neural networks, and present the results of simulations that test the ability of the system on a range of normal and illusory scenes.

## 2 Implementation

Simulations of the model were conducted using the NEXUS Neural Simulator (Sajda and Finkel, 1990; 1991). NEXUS is an interactive simulator designed for modelling multiple interconnected neural maps. The simulator allows considerable flexibility in specifying neuronal properties and neural architectures. The present simulations feature a system of 42 interconnected networks, each of which contains an array of 64x64 units and each of which is topographically organized. Two types of neuronal units are used. Standard neuronal units carry out a linear weighted

summation of their excitatory and inhibitory inputs, and outputs are determined by a sigmoidal function between voltage and firing rate. NEXUS also allows the use of more complex units called PGN (programmable generalized neural) units which execute arbitrary functions or algorithms. A single PGN unit can emulate the function of a small circuit or assembly of standard units.

PGN units are particularly useful in situations in which an intensive computation is being performed but the anatomical and physiological details of the how the operation is performed *in vivo* are unknown. Alternatively, PGN units can be used to carry out functions in a time-efficient manner; for example, to implement a one-step winner-take-all algorithm. The PGN units used in the present simulations can all be replaced with circuits composed of standard neuronal units, but this incurs a dramatic increase in processing time and memory allocation with no change in functional behavior at the system level.

No learning is involved in the network dynamics. The model is intended to correspond to preattentive perception, and the interpretation of even complex scenes requires only a few cycles of network activity. The details of network construction will be described elsewhere; we will focus here on the processes performed and the theoretical issues behind the mechanisms.

### 3 Construction of the Model

The model consists of a number of stages as indicated in Figure 1. The first stage of early visual processing involves networks specialized for the detection of edges, line orientation, line terminations (endstopping), and line junctions (termination of one contour upon another). As Ramachandran (1987) has observed, the visual system must distinguish several different types of edges: we are concerned here with the distinction between edges due to surface discontinuities (transitions between different surfaces) and those due to surface markings (textures, stray lines, etc.). Only the former can be occlusion boundaries. The visual system utilizes several modalities to classify types of edges; we restrict ourselves to a single process carried out by the second processing stage, a network that determines which segments belong to which contours and whether the contours are closed.

When two contours cross each other, forming an "X" junction, there are several possible perceptual interpretations of which arms of the "X" should be joined. Our networks carry out the simple rule that discontinuities should be minimized-i.e., lines and curves should continue

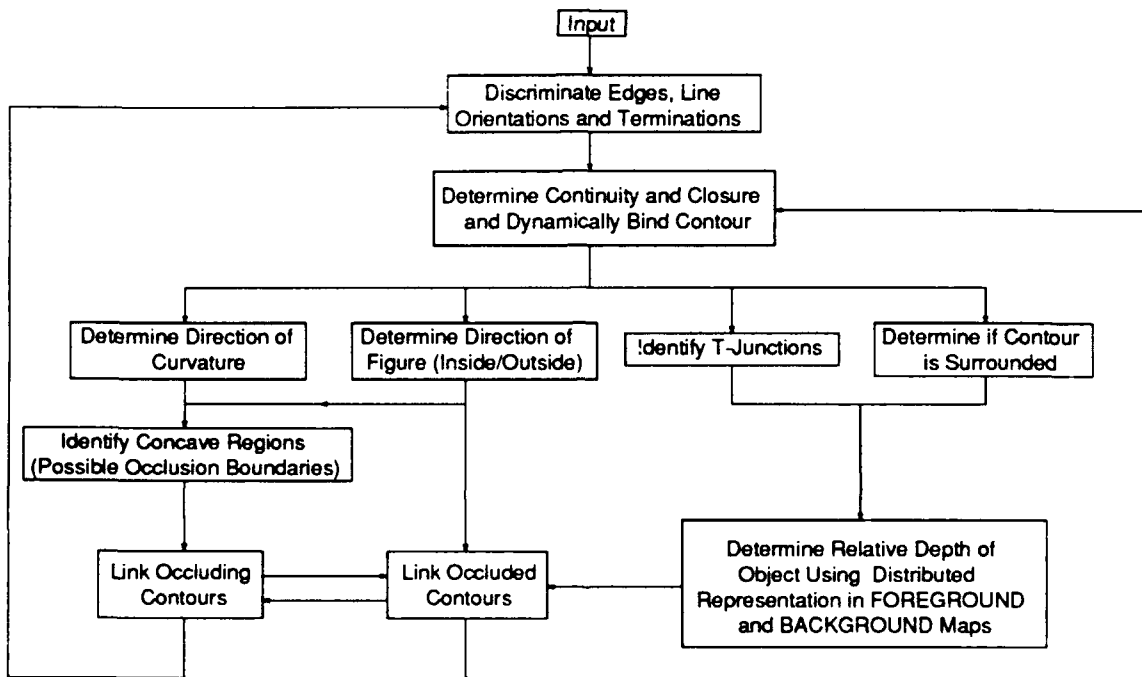


Figure 1: Major processing stages in the model. Each process is carried out by one or more networks. Following early visual stages, information flows through two largely parallel pathways—one concerned with identifying and linking occlusion boundaries (left side) and another concerned with stratifying objects in depth (right side). Networks are multiply interconnected and note the presence of the two reentrant feedback pathways.

as straight (or with as much the same curvature) as possible. Similar assumptions underlie previous models (Ullman, 1976), and this notion is in accord with psychophysical findings that discontinuities contain more information than continuous segments (Attneave, 1954; Resnikoff, 1989). We are thus minimizing the amount of self-generated information.

We employ a simple sequential process to determine whether a contour is closed—each unit on a closed contour requires that at least two of its nearest neighboring units also be on the contour. It is computationally difficult to determine closure in parallel. We speculate that, *in vivo*, the process is carried out by a combination of endstopped units and large-receptive field cells arranged in an architecture similar to that described by Rockland and Lund (1982) in Area 17 (Mitchison and Crick, 1982). Once closure is determined, it is computationally efficient for the units involved to be identified with a “tag”. Several of the higher level processes discussed below require that units responding to the same contour be distinguishable from those responding to different contours. There are several possible physiological mechanisms which could subserve such a tag—one possible mechanism is phase-locked firing (Gray and Singer, 1989; Eckhorn, et al., 1988). We have implemented this contour binding tag through the use of PGN units (section 2), which are capable of representing several distinct tags. It must be emphasized, however, that the model is compatible with a number of possible physiological mechanisms.

Closed contours are a necessary condition to identify a proto-object, but sufficiency requires two additional components. As shown in Figure 1, the remaining determinations are carried out in parallel. One stage is concerned with determining on which side of the contour the figure lies, i.e., distinguishing inside from outside. The problem can be alternatively posed as determining which surface “owns” the contour (Koffka, 1935; Nakayama and Shimojo, 1990). This is a nontrivial problem which, in general, requires global information about the figure. The classic example is the spiral (Minsky and Papert, 1969; Sejnowski and Hinton, 1987) in which it is impossible to determine whether a point is inside or outside based on only local information. The mechanism we employ, as shown in Figure 2, is based upon the following simple observation. Suppose a unit projects its dendrites in a stellate configuration and that the dendrites are activated by units responding to a contour. Then if a given unit is inside a contour, all of its dendrites will be activated (i.e., will intersect the contour); if the unit is outside, then only some of its dendrites will be activated. A winner-take-all interaction between the two units

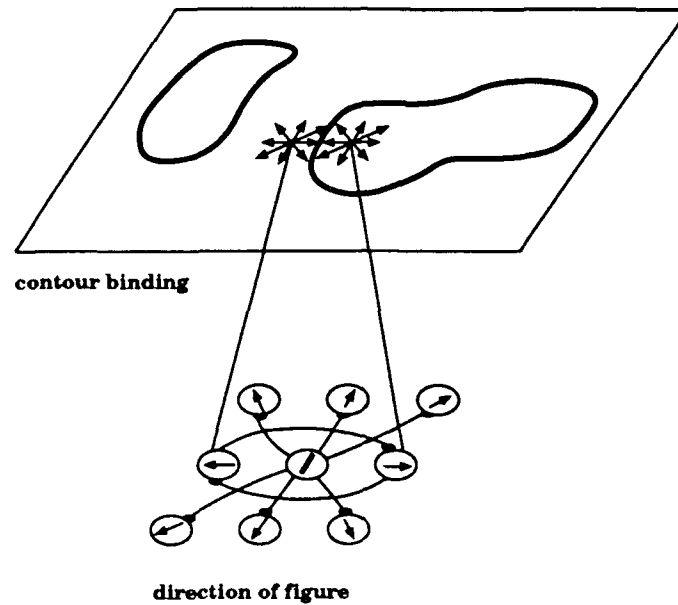


Figure 2: Neural circuit for determining direction of figure (inside vs. outside). Hypothetical input stimulus consists of two closed contours (bold curves). The central unit of 3x3 array (shown below) determines the local orientation of the contour. Surrounding units represent possible directions (indicated by arrows) of the inside of the figure relative to the contour. All surrounding units are inhibited (black circles) except for the two units located perpendicular to local orientation of the contour. Units receive inputs from the contour binding map via dendrites that spread out in a stellate configuration, as indicated by clustered arrows (dendrites extend over long distances in map). Units inside the figure will receive more inputs than those located outside the figure. The two uninhibited units compete in a winner-take-all interaction. Note that inputs from separate objects are not confused due to the tags generated in the contour binding map.

will determine which is more strongly activated, and hence which is inside the figure. As shown in Figure 2, it is advantageous to limit this competition to the two units which are located at positions perpendicular to the local orientation of the contour. As will be shown below (see figures 5-7), this network is quite efficient at locating the interior of figures. It also demonstrates deficiencies similar to those of human perception—for example, it cannot distinguish the inside from the outside of a spiral. The mechanism depends, however, upon the contour binding carried out above. Each unit only considers inputs with the appropriate tag—in this way, inputs from separate contours in the scene are not confused.

Identification of a proto-object also requires that the relative depth of the surface be determined. This is carried out chiefly through the use of T-junctions. As shown in figure 3, a



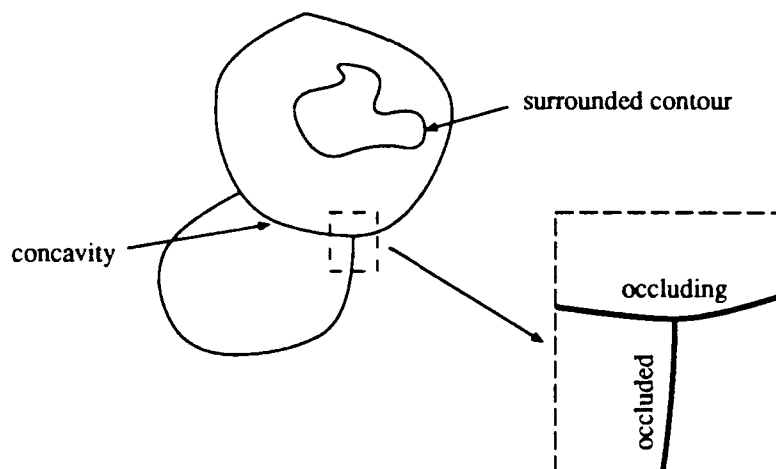


Figure 3: Primary cues for occlusion. T-junctions (shown in the inset) signal a local discontinuity between occluding and occluded contours. Concave regions and surrounded contours suggest occlusion, but are not as reliable indicators as T-junctions. Additional cues such as accretion/deletion of texture (not considered here) are used *in vivo*.

T-junction is formed by the termination of an occluded boundary (the tail of the T) upon an occluding boundary (the head of the T). Although current evidence for cortical units selective for such junctions is lacking, it is trivial to design such a unit based on inputs from orientation selective and endstopped cells.

In this model, T-junctions serve as the major determinant of relative depth. Depth is represented by the relative level of activity in two topographic maps (called *foreground* and *background*). The closest object maximally activates *foreground* units and minimally activates *background* units; the most distant object has the reverse values, and objects located at intermediate depths display intermediate values. The initial state of the two maps is such that all closed contours lie in the background plane. Depth values are then modified at T-junctions—contours forming the head of the “T” are pushed towards the foreground. Since multiple objects can overlap, a contour can be both occluding and occluded—therefore, the relative depth of a contour is determined in a type of push-pull process in which proto-objects are shuffled in depth. The contour binding tag is critical in this process in that all units with the same tag are pushed forward or backward together. (In the more general case of nonplanar objects, the alteration of depth values would depend upon position along the contour)

T-junctions arise in cases of partial occlusion; however, in some instances, a smaller object

may actually lie directly in front of a larger object. In this case, which we call "surround" occlusion, the contour of the occluded object surrounds that of the occluding object. As shown in figure 1, a separate process determines whether such a surround occlusion is present, and in the same manner as T-junctions, leads to a change in the representation of relative depth. The network mechanism for detecting surround occlusion is almost identical to that discussed above for determining the direction of figure (see Figure 2). Note that a similar configuration of two concentric contours arises in the case of a "hole". The model is currently being extended to deal with such non-simply connected objects.

These processes—contour binding, determining direction of the figure, and determination of relative depth—define the proto-object. The remainder of the model is concerned with linking proto-objects into objects. The first step in this endeavor is to identify occlusion boundaries. Since occlusion boundaries are concave segments of contours, such segments must be detected (particularly, concave segments bounded by T-junctions). Although many machine vision algorithms exist for determining convexity, we have chosen to use a simple, neurally plausible mechanism: at each point of a contour, the direction of figure is compared to the direction of curvature (which is determined using endstopped units (Dobbins, et al., 1987)). In convex regions, the two directions are the same; in concave regions, the two directions are opposed. A simple AND mechanism can therefore identify the concave segments of the contours.

Once occlusion borders are identified, proto-objects can be linked by trying to extend, complete, or continue occluded segments. Linkage most commonly occurs between proto-objects in the background, i.e., between spatially separated occluded contours. For example, in figure 3, the occluded contours which terminate at the two T-junctions can be linked to generate a virtual representation of the occluded segment. Since it is impossible to know exactly what the occluded segment looks like, and since it is not actually "perceived", we have chosen not to generate a representation of the occluded segment. Rather, a network link binds together the endpoints of the "tails" of the two T-junctions. In the case where multiple objects are occluded by a single object, the problem of which contours to link can become complex. As shown in figure 4, one important constraint on this process is that the directions of figure be consistent between the two linked proto-objects.

Another condition in which proto-objects can be linked involves the joining of occluding

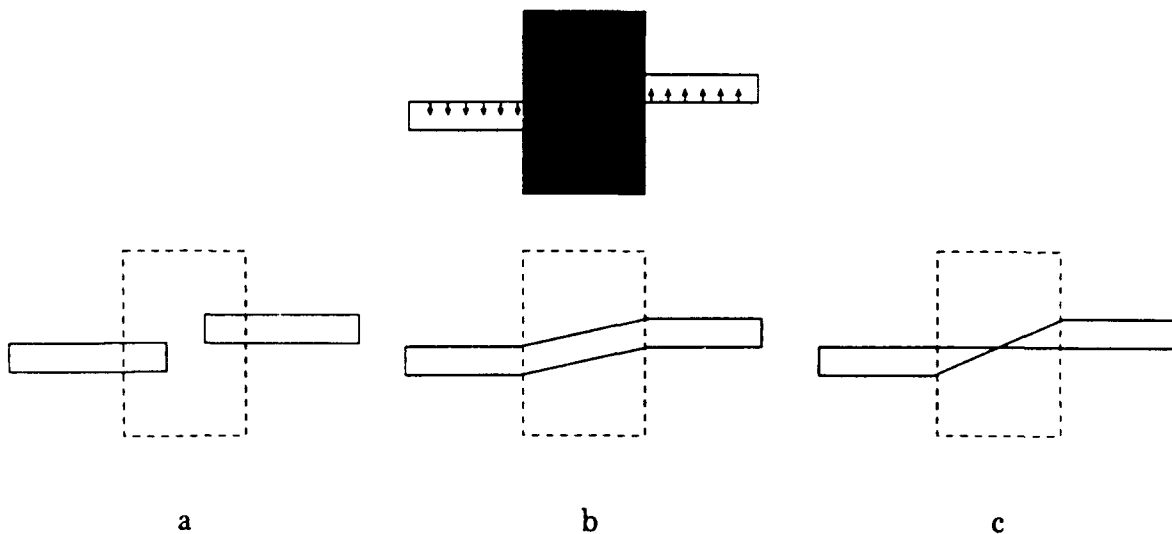


Figure 4: Linking of occluded contours. Three possible perceptual interpretations (below) of an occlusion configuration (above) are shown. Collinearity cannot be the sole criterion for linking occluded edges. Consistency in the direction of figure (inside/outside) between linked objects rules out perception c.

contours, i.e., of proto-objects in the foreground. This phenomenon occurs in our perception of illusory contours, for example, in the Kanizsa triangle (Kanizsa, 1979) or when a gray disc is viewed against a background which changes in a smooth spatial gradient from black to white (Marr, 1982; Shapley and Gordon, 1987). In this case, the heads of the T-junctions are joined, and a representation of the actual contour is generated. The conditions for linkage are that the two contours must be smoothly joined by a line or an arc with constant curvature, and that the direction of figure be consistent (as in the case of occluded contours above).

As indicated in Figure 1, there is an important interaction between these two (occluded and occluding contour) linking processes. Since these links are self-generated by the system (they do not exist in the physical world), they must be scrutinized to avoid false conjunctions. The most powerful check on these processes is their mutual consistency—an increased certainty of the occluded contour continuation being correct increases the confidence of the occluding contour continuation, and vice versa. For example, in the case of the Kanizsa triangle, the “pac-man”-like figures can be completed to form complete circles by simply continuing the contour of the pac-man. The relative ease of completing the occluded contours, in turn, favors the construction of the illusory contours which correspond to the continuations of the occluding contours. In fact,

we believe that the interaction between these two processes determines the perceptual vividness of the illusion.

The final steps in the process involve a recurrent feedback (or reentry, Finkel and Edelman, 1989) from the networks which generate these links back to earlier stages so that the completed contours can be treated as real objects. Note that the occluded contours feedback to the contour binding stage, not to the line discrimination stage, since in this case, the link is virtual, and there is no generated line whose orientation, etc., can be determined. The feedback is particularly important for integrating the outputs of the two parallel paths. For example, once an occluding contour is generated, as in the illusory contours generated in the Kanizsa triangle, it creates a new T-junction (with the circular arc as the "tail" and the illusory contour as the "head" of the "T"). On the next iteration through the system, this T-junction is identified by networks in the other parallel path of the system (see Figure 1), and is used to stratify the illusory contour in depth.

## 4 Results of Simulations

### Linking Proto-Objects

We present the results of three simulations which illustrate the ability of the system to discriminate objects. Figure 5 shows a visual scene that was presented to the system. The early networks discriminate the edges, lines, terminations, and junctions present in the scene. Figure 5A displays the contour binding tags assigned to different scene elements (on the first and fifth cycle of activity). Each box represents elements with a common tag, different boxes represent different tags, and the ordering of the boxes is arbitrary. Note that on the first cycle of activity, discontinuous segments of contours are given separate tags. These tags will later be changed as a result of feedback from the linking processes.

Figure 5B shows the output of the *direction of figure* network, for a small portion of the input scene (near the horse's head). The direction of the arrows indicates the direction of figure determined by the network. The correct direction of figure is determined in all cases: for the horse's head, and for the horizontal and vertical posts of the fence. Once the direction of figure is identified, occluded contours can be linked (as in figure 4), and proto-objects combined into objects. This linkage is what changes the contour binding tags, so that after several cycles

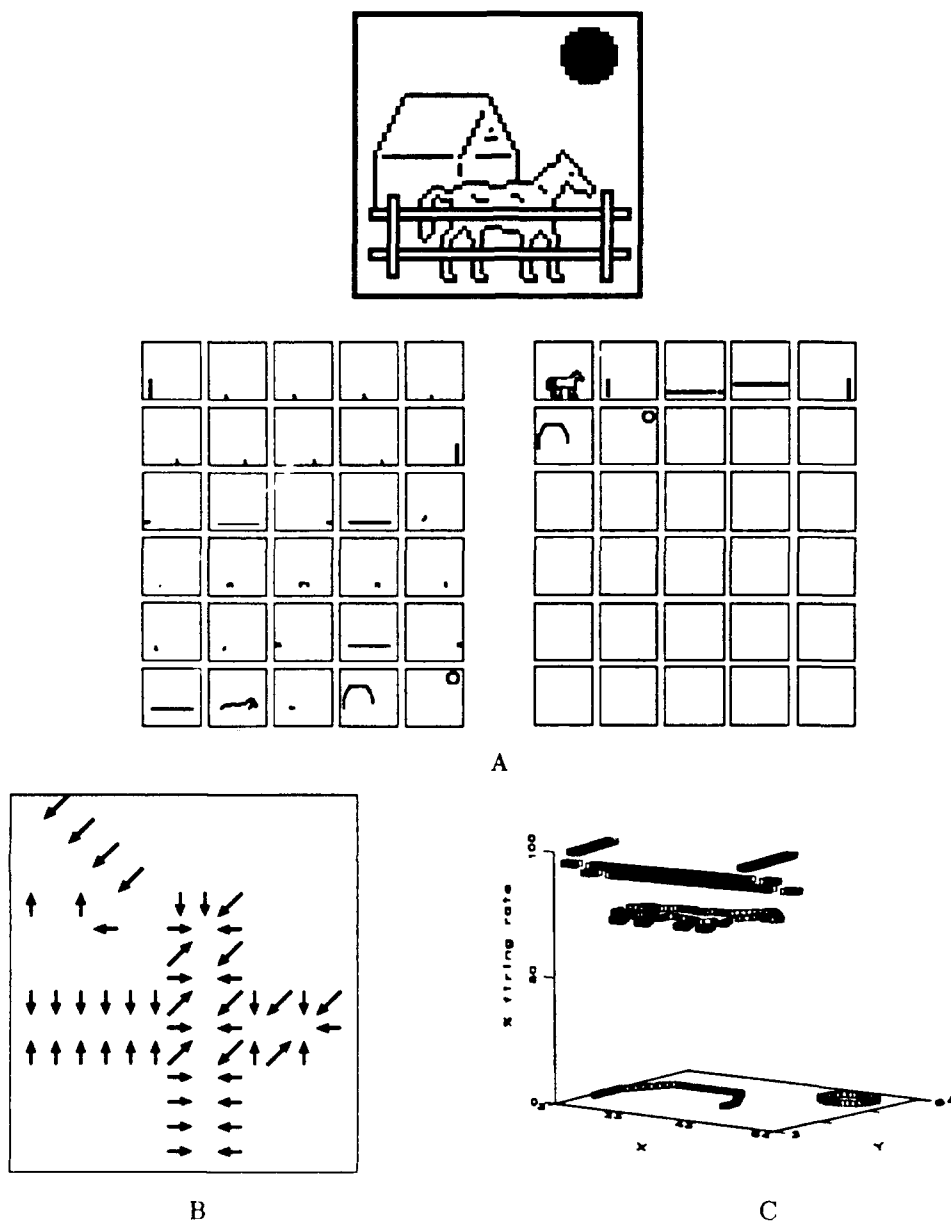


Figure 5: Object discrimination and stratification in depth. Top panel shows a 64 x 64 input stimulus presented to the system. **A** Spatial histogram of the contour binding tags (each box shows units with common tag, different boxes represent different tags, and the order of the boxes is arbitrary). Initial tags shown on left; tags after five iterations shown on right. Note that linking of occluded contours has transformed proto-objects into objects. **B** Magnified view of a local section of the direction of figure network corresponding to portion of the image near horse's nose and crossing fence posts. Arrows indicate direction of inside of proto-objects as determined by network. **C** Relative depth of objects in scene as determined by the system. Plot of activity (% of maximum) of units in the foreground network after 5 iterations. Points with higher activity are "perceived" as being relatively closer to the viewer.

(Figure 5A, right), separate tags are assigned to separate objects—the horse, the gate posts, the house, the sun.

The presence of T-junctions (e.g., between the horse's contour and the fence, between the house and the horse's back) are used by the system to force various objects into different depth planes. The results of this process are displayed in Figure 5C which plots the firing rate (percent of maximum) of units in the *foreground* network. The system has successfully stratified the fence, horse, house and sun. The actual depth value determined for each object is somewhat arbitrary, and can vary depending upon minor changes in the scene—the system is designed only to achieve the correct relative ordering, not absolute depth. Note that the horizontal and vertical posts of the fence are perceived at different depths—this is because of the T-junctions present between them; in fact, the two surfaces do lie at slightly different depths. In addition, there is no way to determine the relative depth of the two objects in the background, the house and the sun, because they bear no occlusion relationship to each other. Again, this conforms to human perceptions, e.g., the sun and the moon appear about the same distance away. The system thus appears to process occlusion information in a manner similar to human perception.

### **Gestalt Psychology of a Network**

The system also displays a response consistent with human responses to a number of illusory stimuli. Figure 6 shows a stimulus, adapted from an example of Kanizsa (1979), which shows that preservation of local continuity in contours is more powerful than global symmetry in perception (this is contrary to classical Gestalt theory—e.g., Koffka, 1935). As shown in the middle panels, there are two possible perceptual interpretations of the contours—on the left, the two figures respect local continuity (this is the dominant human perception); on the right, the figures respect global symmetry.

Figure 6A shows the contour binding tags assigned by the system to this stimulus, and figure 6B shows the direction of figure that was determined. Both results indicate that the network makes the same perceptual interpretation as a human observer.

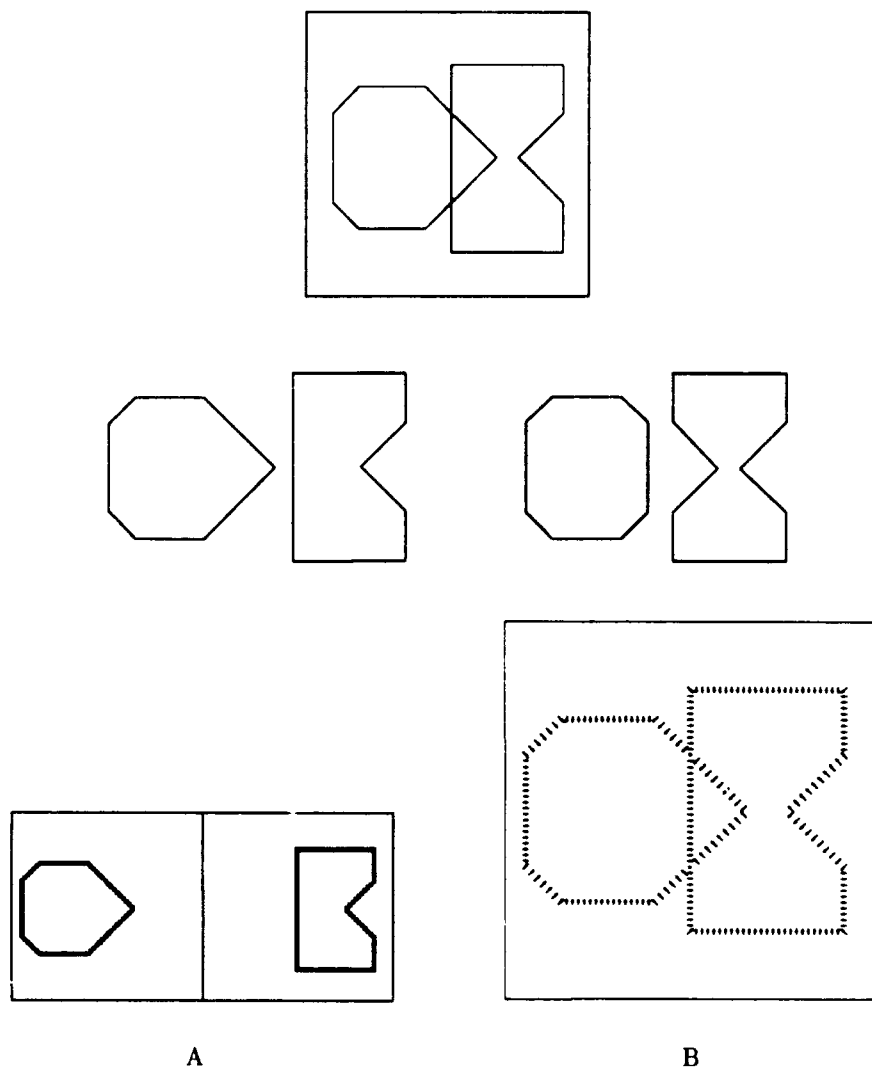


Figure 6: Minimization of ambiguous discontinuities. Upper panel shows an ambiguous stimulus (adapted from Kanizsa, 1979), two possible perceptual interpretations of which are shown below. The interpretation on the left is dominant for humans, despite the figural symmetry of the segmentation on the right. Stimulus was presented to the system, results shown after three iterations. **A** Spatial histogram showing the contour binding patterns (as in 5B). The network segments the figures in the same manner as human perception. **B** Determination of direction of figure confirms network interpretation (note at junction points, direction of figure is indeterminate).

## Occlusion Capture

The final simulation shows the ability of the system to generate illusory contours and to use illusory objects in a veridical fashion. The stimulus is, again, adapted from Kanizsa (1979), and shows a perceptually vivid, illusory white square in a field of black discs. The illusory square appears to be closer to the viewer than the background, and, in addition, the four discs which lie inside its borders appear even closer than the square. This is an example of what we call "occlusion capture", an effect related to Ramachandran's capture phenomenon (Ramachandran and Cavanaugh, 1985; Ramachandran, 1986), in which the illusory square has "captured" the discs within its borders and they are thus pulled into the foreground.

Figure 7A shows the contour binding tags after one (left) and three (right) cycles of activity. Each disc receives a separate tag. After the responses to illusory square are generated, the illusory contours are fed back to the contour binding network and given a common tag. Note that the edges of the discs occluded by the illusory square are now given the same tag as the square, not the same tags as the discs.

The change in "ownership" of the occluded edges of the discs is the critical step in defining the illusory square as an object. For example, Figure 7B shows the output of the *direction of figure* network after one and three cycles of activity. The large display shows that every disc is identified as an object with the inside of the disc correctly labeled in each case. The two insets focus on a portion of the display near the bottom left edge of the illusory square. At first, the system identifies the "L"-shaped angular edge as belonging to the disc, and thus the direction of figure arrows point "inward". After three cycles of activity, this same "L"-shaped edge is identified as belonging to the illusory square, and thus the arrows now point towards the inside of the square, rather than the inside of the disc. This change in the ownership of the edge results from the discrimination of occlusion—the edge has been determined to be an occlusion border. The interconnected processing of the system then results in a change in the direction of figure and of the continuity tags associated with this edge. The illusory square is perceived as an *object*. Its four contours are bound together, the contours are bound to the internal surface, and the properties of the surface are identified.

Figure 7C displays the firing rate of units in the *foreground* map (as in 5C), thus showing the



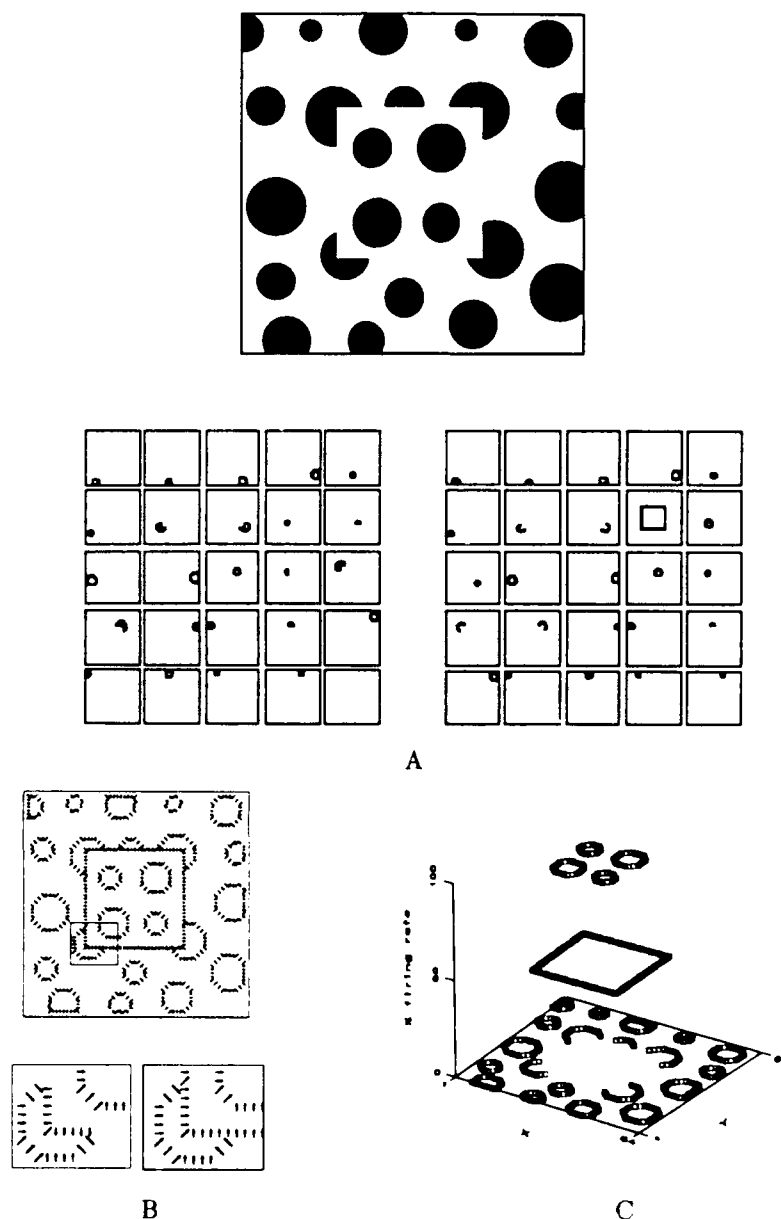


Figure 7: Occlusion capture. Upper panel shows stimulus (adapted from Kanizsa, 1979) in which we perceive a white illusory square. Note that the four black discs inside the illusory square appear closer than the background. A 64 x 64 discrete version of stimulus was presented to the network. **A** Spatial histogram (as in 5A) of the initial and final (after 3 iterations) contour binding tags. Note that the illusory square is bound as an object. **B** Direction of figure determined by the system. Insets show a magnified view of the initial (left) and final (right) direction of figure (region of magnification is indicated). Note that the direction of figure of the "mouth" of the pac-man flips once the illusory contour is generated. **C** Activity in the *foreground* network (% of maximum) demonstrates network stratification of objects in relative depth. The illusory square has "captured" the background texture.

relative depths discriminated by the system. The discs are placed in the background, the illusory square at an intermediate depth, and the discs located within the borders of the illusory square are located closest to the viewer. In this case, the depth cue which forces the internal discs to the foreground is not due to T-junctions, but rather to surround occlusion (see Figure 3). Once the illusory square is generated, the contours of the discs inside the square are surrounded by that of the square. The fact that the contour is "illusory" is irrelevant, since once responses are generated in the networks responsible for linking occluding contours and are then fed back to earlier networks, they are indistinguishable from responses to real contours in the periphery. Thus the system demonstrates occlusion capture corresponding to human perceptions of this stimulus.

## 5 Discussion

In most visual scenes, the majority of objects are partially occluded. Our seamless perception of the world depends upon an ability to complete or link the spatially separated, non-occluded portions of an object. We have used the idea that the visual system identifies proto-objects (which may or may not be objects) and then attempts to link these proto-objects into larger structures. This linking process is most apparent in the perception of illusory contours, and our model can account for a wide range of these illusions. We have only considered static visual scenes, however, one of the major cues to the linking process is common motion of proto-objects. During development, common motion may, in fact, play the largest role in establishing our concept of what is an object (Terminé, et al., 1987). Object definition also clearly depends upon higher cognitive processes such as attention, context and categorization (Rosch and Lloyd, 1978), and such processes must eventually be considered.

This model builds upon previous neural, psychological, and machine vision studies. Several models of illusory contour generation (Ullman, 1976; Peterhans and von der Heydt, 1989; Finkel and Edelman, 1989) have used related mechanisms to check for collinearity and to generate the illusory contours. Our model differs at a more fundamental level—we are concerned with objects not just contours. To define an object, surfaces must also be considered. For example, in a simple line drawing, we perceive an interior surface despite the fact that no surface properties are indicated. Thus, the model must be capable of characterizing a surface—and it does so, in a

rudimentary manner, by determining the direction of figure and relative depth. Nakayama and Shimojo (1990) have approached the problem of surface representation from a similar viewpoint. They discuss how contours and surfaces become associated, how T-junctions serve to stratify objects in depth, and how occluded surfaces are amodally completed. However, Nakayama eschews the construction of a bottom-up model and instead explores the external "ecological" constraints upon perception. In addition to these Gibsonian constraints, we emphasize the importance of *internal* constraints imposed by physiological mechanisms and neural architectures. Nakayama has also explored the interactions between occlusion and surface attributes. A more complete model must consider such surface properties such as color, brightness, texture, and surface orientation. The examination of how surface features might interact with contour boundaries has been pioneered by Grossberg (1987). Finally, in some regards, our model constitutes the first step of a "bottom-up" model of object perception (Kanizsa, 1979; Biederman, 1987). It is interesting that regardless of one's orientation (bottom-up or top-down) the constraints of the physical problem result in certain similarities of solution as witnessed by the analogies present with AI based models (Fisher, 1989).

One of the most speculative aspects of the model is the use of tags to identify elements as belonging to the same object. Tags linking units responding to the same contour are used to determine the direction of figure and to change the perceived depth of the entire contour based on occlusion relationships detected at isolated points (the T-junctions). It is possible to derive alternative mechanisms for these processes which do not depend upon the use of tags, but they are conceptually inelegant and computationally unwieldy. Our model offers no insight as to the biophysical basis of such a tag, or even whether the tag is implemented in a mechanism based on common time, phase, frequency, or map position. However, the model does place constraints on the spatial and temporal properties of the mechanism. For example, suppose that a phase-dependent mechanism were used with voltage oscillations at 50 Hz (Konig and Schillen, 1991). Assuming that neurons require at least 2 milliseconds to process each tag, each neuron could respond to a maximum of 10 tags (2-4 tags is probably a more reasonable estimate). This number should correspond to the number of objects that can be simultaneously discriminated. One could thus imagine a mechanism in which only a handful of objects can be concurrently attended to, and units responding to each object fire in a distinct time window (e.g., the first 10

milliseconds of every 50 millisecond interval).

At the outset, we discussed the importance of both binding and segmentation for visual object discrimination. Our model has largely dealt with the segmentation problem, however, the two problems are not entirely independent. For example, the association of a depth value with the object discriminated is, in essence, an example of the binding of an attribute to an object. Consideration of additional attributes makes the problem more complex, but it also aids in the discrimination of separate objects (Damasio, 1989; Crick and Koch, 1990).

As Ullman (1989) has pointed out, it is not *logically* necessary for object discrimination to take place before object recognition can occur. However, if one considers the function of the multiple extrastriate visual areas leading to inferotemporal cortex, it appears reasonable that the visual system is using all the processes at its disposal to generate meaningful representations of the visual scene. The question of whether you have to know that something is a "thing" before you can recognize what kind of thing it is, remains to be determined through psychophysical experiment.

Nonetheless, our model shows that one can build a self-contained system for discriminating objects based on occlusion relationships. The model is successful at stratifying simple visual scenes, for linking the representations of occluded objects, and at generating responses to illusory objects in a manner consistent with human perceptual responses. The model uses neural circuits that are biologically-based, and conforms to general neural principles, such as the use of a distributed representation for depth. The system can be tested in psychophysical paradigms and the results compared to human and animal results. In this manner, a computational model which is designed based on physiological data and tested with psychophysical data offers a powerful paradigm for bridging the gap between neuroscience and perception.

## Acknowledgements

This work was supported by grants from The Office of Naval Research (N0014-90-J-1864), The Whitaker Foundation, and The McDonnell-Pew Program in Cognitive Neuroscience.

## References

- [1] J. Aloimonos and D. Shulman. *Integration of Visual Modules*. Academic Press, New York,

- 1989.
- [2] F. Attneave. Some informational aspects of visual perception. *Psychology Review*, 61:183-193, 1954.
  - [3] H. B. Barlow. Critical limiting factors in the design of the eye and visual cortex. *Proc. Royal Society (London)*, B212:1-34, 1981.
  - [4] I. Biederman. Recognition by components: a theory of human image understanding. *Psychological Review*, 94:115-147, 1987.
  - [5] F. Crick and C. Koch. Towards a neurobiological theory of consciousness. *Seminars in Neuroscience*, 2:263-275, 1990.
  - [6] A. R. Damasio. The brain binds entities and events by multiregional activation from convergence zones. *Neural Computation*, 1:1223-1232, 1989.
  - [7] A. S. Dobbins, S. W. Zucker, and M. S. Cynader. Endstopping in the visual cortex as a neural substrate for calculating curvature. *Nature*, 329:438-441, 1987.
  - [8] R. Eckhorn, R. Bauer, W. Jordan, M. Brosch, W. Kruse, M. Munk, and H. Reitboeck. Coherent oscillations: a mechanism of feature linking in the visual cortex? *Biological Cybernetics*, 60:121-130, 1988.
  - [9] L. Finkel and G. Edelman. Integration of distributed cortical systems by reentry: a computer simulation of interactive functionally segregated visual areas. *Journal of Neuroscience*, 9:3188-3208, 1989.
  - [10] R. B. Fisher. *From Objects to Surfaces*. John Wiley & Sons, New York, 1989.
  - [11] C. M. Gray and W. Singer. Neuronal oscillations in orientation columns of cat visual cortex. *Proceedings of the National Academy of Science USA*, 86:1698-1702, 1989.
  - [12] R. L. Gregory. Cognitive contours. *Nature*, 238:51-52, 1972.
  - [13] S. Grossberg. Cortical dynamics of three-dimensional form, color, and brightness perception, i: monocular theory. *Perception and Psychophysics*, 41:87-116, 1987.

- [14] S. Grossberg and E. Mingolla. Neural dynamics of form perception: boundary completion, illusory figures, and neon color spreading. *Psychology Review*, 92:173-211, 1985.
- [15] A. Guzman. Decomposition of a visual scene into three-dimensional bodies. *Fall Joint Computer Conference*, 1968:291-304, 1968.
- [16] G. Kanizsa. *Organization in Vision*. Praeger, New York, 1979.
- [17] K. Koffka. *Principles of Gestalt Psychology*. Harcourt, Brace, New York, 1935.
- [18] P. Konig and T. Schillen. Stimulus-dependent assembly formation of oscillatory responses: i. synchronization. *Neural Computation*, 3:155-166, 1991.
- [19] S. Lehky and T. Sejnowski. Neural model of stereoacuity and depth interpolation based on distributed representation of stereo disparity. *Journal of Neuroscience*, 7:2281-2299, 1990.
- [20] D. Marr. *Vision: A computational investigation into the Human Representation and Processing of Visual Information*. W. H. Freeman, San Francisco, 1982.
- [21] M. Minsky and S. Papert. *Perceptrons*. MIT Press, Cambridge, MA, 1969.
- [22] G. Mitchison and F. Crick. Long axons within the striate cortex: their distribution, orientation, and patterns of connections. *Proceedings of the National Academy of Science*, 79:3661-3665, 1982.
- [23] K. Nakayama and S. Shimojo. Toward a neural understanding of visual surface representation. *Cold Spring Harbor Symposia on Quantitative Biology*, LV:911-924, 1990.
- [24] E. Peterhans and R. von der Heydt. Mechanisms of contour perception in monkey visual cortex. ii. contours bridging gaps. *Journal of Neuroscience*, 9:1749-1763, 1989.
- [25] G. F. Poggio, F. Gonzalez, and F. Krause. Stereoscopic mechanisms in monkey visual cortex: binocular correlation and disparity selectivity. *Journal of Neuroscience*, 8:4531-4550, 1988.
- [26] T. Poggio, E. B. Gamble, and J. J. Little. Parallel integration of vision modules. *Science*, 242:436-440, 1988.

- [27] V. S. Ramachandran. Capture of stereopsis and apparent motion by illusory contours. *Perception and Psychophysics*, 39:361-373, 1986.
- [28] V. S. Ramachandran. Visual perception of surfaces: a biological theory. In S. Petry and G. E. Meyer, editors, *The Perception of Illusory Contours*, pages 93-108, Springer-Verlag, New York, 1987.
- [29] V. S. Ramachandran and P. Cavanagh. Subjective contours capture stereopsis. *Nature*, 317:527-530, 1985.
- [30] H. L. Resnikoff. *The Illusion of Reality*. Springer-Verlag, New York, 1989.
- [31] K. S. Rockland and J. S. Lund. Widespread periodic intrinsic connections in the tree shrew visual cortex. *Science*, 215:1532-1534, 1982.
- [32] E. Rosch and B. B. Lloyd. *Cognition and Categorization*. Lawrence Erlbaum Associates, Hillsdale, N.J., 1978.
- [33] A. Rosenfeld. Computer vision. *Advances in Computers*, 27:265-308, 1988.
- [34] P. Sajda and L. Finkel. NEXUS: A simulation environment for large-scale neural systems. *SIMULATION*, submitted.
- [35] P. Sajda and L. Finkel. The NEXUS neural simulation environment. *University of Pennsylvania Technical Report*, 1990.
- [36] T. Sejnowski and G. Hinton. Separating figure from ground with a Boltzmann machine. In M. Arbib and A. Hanson, editors, *Vision, Brain and Cooperative Computation*, pages 703-724, MIT Press, Cambridge, MA, 1987.
- [37] R. Shapley and J. Gordon. The existence of interpolated illusory contours depends on contrast and spatial separation. In S. Petry and G. E. Meyer, editors, *The Perception of Illusory Contours*, pages 109-115, Springer-Verlag, New York, 1987.
- [38] N. Termine, T. Hrynicky, R. Kestenbaum, H. Gleitman, and E. S. Spelke. Perceptual completion of surfaces in infancy. *Journal Experimental Psychology-Human Perception*, 13:524-532, 1987.